# Exploratory data analysis (EDA) using python: a tutorial

**Presentation** · October 2024

DOI: 10.13140/RG.2.2.20542.09284/1

1 author:

Sami D. Alaruri سامي العاروري
Independnet Researcher
79 PUBLICATIONS   530 CITATIONS

SEE PROFILE

# AGENDA

- **Introduction**

- **Overview**

- **Data Inspection & Cleaning Steps**

- **Graphical EDA**

- **Conclusions**

- **References**

# INTRODUCTION: WHAT IS EDA?

Exploratory data analysis is an analysis technique to analyze and investigate the data set and summarize the main characteristics of the dataset. Main advantage of EDA is providing the data visualization of data after conducting the analysis.

Tukey defined data analysis in 1961 as: "Procedures for analyzing data, techniques for interpreting the results of such procedures, ways of planning the gathering of data to make its analysis easier, more precise or more accurate, and all the machinery and results of (mathematical) statistics which apply to analyzing data".
https://en.wikipedia.org/wiki/Exploratory_data_analysis

# OVERVIEW

Obesity is a complex disease involving having too much body fat. Obesity isn't just a cosmetic concern. It's a medical problem that increases the risk of many other diseases and health problems. These can include heart disease, diabetes, high blood pressure, high cholesterol, liver disease, sleep apnea and certain cancers.

There are many reasons why some people have trouble losing weight. Often, obesity results from inherited, physiological and environmental factors, combined with diet, physical activity and exercise choices.

The good news is that even modest weight loss can improve or prevent the health problems associated with obesity. A healthier diet, increased physical activity and behavior changes can help you lose weight. Prescription medicines and weight-loss procedures are other options for treating obesity. In this Exploratory Data Analysis (EDA) using Python we will examine the NObeyesdad (Obesity level of the individual) as a function of several factors (i.e., age, gender, height, alcohol consumption ..etc.).

https://www.mayoclinic.org/diseases-conditions/obesity/symptoms-causes/syc-20375742

# DATASET INFORMATION

This dataset include data for the estimation of obesity levels in individuals from the countries of Mexico, Peru and Colombia, based on their eating habits and physical condition. The data contains 17 attributes and 2111 records, the records are labeled with the class variable NObesity (Obesity Level), that allows classification of the data using the values of Insufficient Weight, Normal Weight, Overweight Level I, Overweight Level II, Obesity Type I, Obesity Type II and Obesity Type III. 77% of the data was generated synthetically using the Weka tool and the SMOTE filter, 23% of the data was collected directly from users through a web platform.

https://archive.ics.uci.edu/dataset/544/estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition

# DATA INSPECTION & CLEANING STEPS

1. Dataset dimensions

2. Titles of columns

3. Data types

4. Missing values

5. Nulls in the dataset (not required for this dataset)

6. Duplicate rows

7. NaN values

8. Infinity values

9. Outliers detection

10. Encode categorical features

# LOADING PYTHON LIBRARIES

```python
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import pandas as pd
import numpy as np
from sklearn.preprocessing import OneHotEncoder, StandardScaler,LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
warnings.filterwarnings("ignore")
```

**pandas (https://pandas.pydata.org)**

**numpy (https://numpy.org)**

**matplotlib (https://matplotlib.org)**

**seaborn (https://seaborn.pydata.org)**

**sklearn (https://scikit-learn.org/stable)**

# DATA SET ATTRIBUTES

1. Age: The age of the individual.
2. Gender: The gender of the individual (e.g., Male, Female).
3. Height: The height of the individual in meters.
4. Weight: The weight of the individual in kilograms.
5. CALC: Unknown column. You might need to check the data source or documentation to understand what this column represents.
6. FAVC: Whether the individual frequently consumes high caloric food (e.g., yes, no).
7. FCVC: Frequency of consumption of vegetables (numeric scale).

8. NCP: Number of main meals per day (numeric scale).
9. SCC: Squamous cell carcinoma
10. SMOKE: Whether the individual smokes (e.g., yes, no).
11. CH2O: Consumption of water daily (numeric scale).
12. family_history_with_overweight: Whether the individual has a family history of overweight (e.g., yes, no).
13. FAF: Physical activity frequency (numeric scale).
14. TUE: Time using technology devices (numeric scale).
15. CAEC: Unknown column. You might need to check the data source or documentation to understand what this column represents.
16. MTRANS: Mode of transportation (e.g., Public Transportation, Walking).
17. NObeyesdad: Obesity level of the individual (e.g., Normal_Weight, Overweight_Level_I).
    **Note**: The data set contains a mixture of numeric and categorical data.

# PANDAS FUNCTIONS FOR DATA INSPECTION

- df.head()/df.tail()

- df.sample()

- df.info()

- df.columns

- df.describe()

- df.shape

- Df.count()

https://www.geeksforgeeks.org/pandas-functions-in-python/

# USING PANDAS FOR LOADING THE DATA FILE & VIEWING THE FIRST 5 ROWS



```
data = pd.read_csv("/content/sample_data/ObesityDataSet_raw_and_data.csv") # Load data set using pandas
data.head(5)
```

| | Age | Gender | Height | Weight | CALC | FAVC | FCVC | NCP | SCC | SMOKE | CH2O | family_history_with_overweight | FAF | TUE | ( |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 21.0 | Female | 1.62 | 64.0 | no | no | 2.0 | 3.0 | no | no | 2.0 | yes | 0.0 | 1.0 | Some |
| 1 | 21.0 | Female | 1.52 | 56.0 | Sometimes | no | 3.0 | 3.0 | yes | yes | 3.0 | yes | 3.0 | 0.0 | Some |
| 2 | 23.0 | Male | 1.80 | 77.0 | Frequently | no | 2.0 | 3.0 | no | no | 2.0 | yes | 2.0 | 1.0 | Some |
| 3 | 27.0 | Male | 1.80 | 87.0 | Frequently | no | 3.0 | 3.0 | no | no | 2.0 | no | 2.0 | 0.0 | Some |
| 4 | 22.0 | Male | 1.78 | 89.8 | Sometimes | no | 2.0 | 1.0 | no | no | 2.0 | no | 0.0 | 0.0 | Some |

**Dataset source: UC Irvine-Machine Learning Repository**
https://archive.ics.uci.edu/dataset/544/estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition

# EXAMINING 5 ROWS IN RANDOM & AT THE END OF THE DATASET

```
1 data.sample(5)
```

| | Gender | Age | Height | Weight | family_history_with_overweight | FAVC | FCVC | NCP | CAEC | SMOKE | CH2O | SCC | FAF | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1679** | Male | 31.761799 | 1.751688 | 119.205308 | yes | yes | 2.149610 | 3.000000 | Sometimes | no | 2.133876 | no | 0.393452 | 0.3 |
| **752** | Male | 21.142432 | 1.855353 | 86.413388 | yes | yes | 2.000000 | 3.000000 | Sometimes | no | 1.345298 | no | 1.097905 | 1.0 |
| **1546** | Male | 25.298400 | 1.827279 | 120.996074 | yes | yes | 3.000000 | 3.000000 | Sometimes | no | 3.000000 | no | 1.110215 | 0.3 |
| **1741** | Male | 28.255199 | 1.816547 | 120.699119 | yes | yes | 2.997951 | 3.000000 | Sometimes | no | 2.715856 | no | 0.739881 | 0.9 |
| **1790** | Male | 23.147644 | 1.815514 | 120.337664 | yes | yes | 2.996717 | 2.791366 | Sometimes | no | 2.626309 | no | 1.194898 | 0.0 |

```
1 data.tail(5)
```

| | Gender | Age | Height | Weight | family_history_with_overweight | FAVC | FCVC | NCP | CAEC | SMOKE | CH2O | SCC | FAF | TUE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **2106** | Female | 20.976842 | 1.710730 | 131.408528 | yes | yes | 3.0 | 3.0 | Sometimes | no | 1.728139 | no | 1.676269 | 0.906247 |
| **2107** | Female | 21.982942 | 1.748584 | 133.742943 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.005130 | no | 1.341390 | 0.599270 |
| **2108** | Female | 22.524036 | 1.752206 | 133.689352 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.054193 | no | 1.414209 | 0.646288 |
| **2109** | Female | 24.361936 | 1.739450 | 133.346641 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.852339 | no | 1.139107 | 0.586035 |
| **2110** | Female | 23.664709 | 1.738836 | 133.472641 | yes | yes | 3.0 | 3.0 | Sometimes | no | 2.863513 | no | 1.026452 | 0.714137 |

11

# USING PANDA FUNCTIONS FOR VIEWING THE DATA



**Data set size:**
**2111 rows x 17 columns**

# DATA STATISTICS & MISSING VALUES CHECK



**Dataset Statistics Summary**



**Missing Values Check**
**No missing values**

# INSPECTING THE DATASET FOR DUPLICATE ROWS & DROPPING THE DUPLICATE ROWS

```python
# Check for duplicate rows in the data set
duplicate_rows = data.duplicated().sum()
print("Duplicate Rows:")
print(duplicate_rows)
```

```
Duplicate Rows:
24
```

**As shown above, there is 24 duplicate rows in the data set.**

```python
data=data.drop_duplicates()  # use dataframe.drop_duplicates() to drop the duplicate rows
display(data)
```

| | Age | Gender | Height | Weight | CALC | FAVC | FCVC | NCP | SCC | SMOKE | CH2O | family_history_with_overweight |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 21.000000 | Female | 1.620000 | 64.000000 | no | no | 2.0 | 3.0 | no | no | 2.000000 | yes |
| 1 | 21.000000 | Female | 1.520000 | 56.000000 | Sometimes | no | 3.0 | 3.0 | yes | yes | 3.000000 | yes |
| 2 | 23.000000 | Male | 1.800000 | 77.000000 | Frequently | no | 2.0 | 3.0 | no | no | 2.000000 | yes |
| 3 | 27.000000 | Male | 1.800000 | 87.000000 | Frequently | no | 3.0 | 3.0 | no | no | 2.000000 | no |
| 4 | 22.000000 | Male | 1.780000 | 89.800000 | Sometimes | no | 2.0 | 1.0 | no | no | 2.000000 | no |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 2106 | 20.976842 | Female | 1.710730 | 131.408528 | Sometimes | yes | 3.0 | 3.0 | no | no | 1.728139 | yes |
| 2107 | 21.982942 | Female | 1.748584 | 133.742943 | Sometimes | yes | 3.0 | 3.0 | no | no | 2.005130 | yes |
| 2108 | 22.524036 | Female | 1.752206 | 133.689352 | Sometimes | yes | 3.0 | 3.0 | no | no | 2.054193 | yes |
| 2109 | 24.361936 | Female | 1.739450 | 133.346641 | Sometimes | yes | 3.0 | 3.0 | no | no | 2.852339 | yes |
| 2110 | 23.664709 | Female | 1.738836 | 133.472641 | Sometimes | yes | 3.0 | 3.0 | no | no | 2.863513 | yes |

2087 rows × 17 columns

```python
data.shape # check data set dimension after removing duplicates
```
```
(2087, 17)
```

**24 duplicate row**
**After dropping the 24 duplicate Rows, new dataset size:**
**2087 rows x 17 columns**

# NAN VALUES CHECK



**No NaN values in the dataset**

# INFINITY VALUES CHECK

```python
# Check the presence of infinity values in the data set
data.isin([np.inf, -np.inf])
print(data)
```

No ±∞ values in the data set

```
        Age  Gender    Height      Weight         CALC FAVC  FCVC  NCP  \
0    21.000000  Female  1.620000   64.000000           no   no   2.0  3.0
1    21.000000  Female  1.520000   56.000000    Sometimes   no   3.0  3.0
2    23.000000    Male  1.800000   77.000000   Frequently   no   2.0  3.0
3    27.000000    Male  1.800000   87.000000   Frequently   no   3.0  3.0
4    22.000000    Male  1.780000   89.800000    Sometimes   no   2.0  1.0
...        ...     ...       ...         ...          ...  ...   ...  ...
2106 20.976842  Female  1.710730  131.408528    Sometimes  yes   3.0  3.0
2107 21.982942  Female  1.748584  133.742943    Sometimes  yes   3.0  3.0
2108 22.524036  Female  1.752206  133.689352    Sometimes  yes   3.0  3.0
2109 24.361936  Female  1.739450  133.346641    Sometimes  yes   3.0  3.0
2110 23.664709  Female  1.738836  133.472641    Sometimes  yes   3.0  3.0

      SCC SMOKE      CH2O family_history_with_overweight         FAF       TUE  \
0      no    no  2.000000                            yes  0.000000  1.000000
1     yes   yes  3.000000                            yes  3.000000  0.000000
2      no    no  2.000000                            yes  2.000000  1.000000
3      no    no  2.000000                             no  2.000000  0.000000
4      no    no  2.000000                             no  0.000000  0.000000
...   ...   ...       ...                            ...       ...       ...
2106   no    no  1.728139                            yes  1.676269  0.906247
2107   no    no  2.005130                            yes  1.341390  0.599270
2108   no    no  2.054193                            yes  1.414209  0.646288
2109   no    no  2.852339                            yes  1.139107  0.586035
2110   no    no  2.863513                            yes  1.026452  0.714137

           CAEC                 MTRANS            NObeyesdad
0     Sometimes  Public_Transportation         Normal_Weight
1     Sometimes  Public_Transportation         Normal_Weight
2     Sometimes  Public_Transportation         Normal_Weight
3     Sometimes                Walking     Overweight_Level_I
4     Sometimes  Public_Transportation    Overweight_Level_II
...         ...                    ...                    ...
2106  Sometimes  Public_Transportation       Obesity_Type_III
2107  Sometimes  Public_Transportation       Obesity_Type_III
2108  Sometimes  Public_Transportation       Obesity_Type_III
2109  Sometimes  Public_Transportation       Obesity_Type_III
2110  Sometimes  Public_Transportation       Obesity_Type_III

[2087 rows x 17 columns]
```
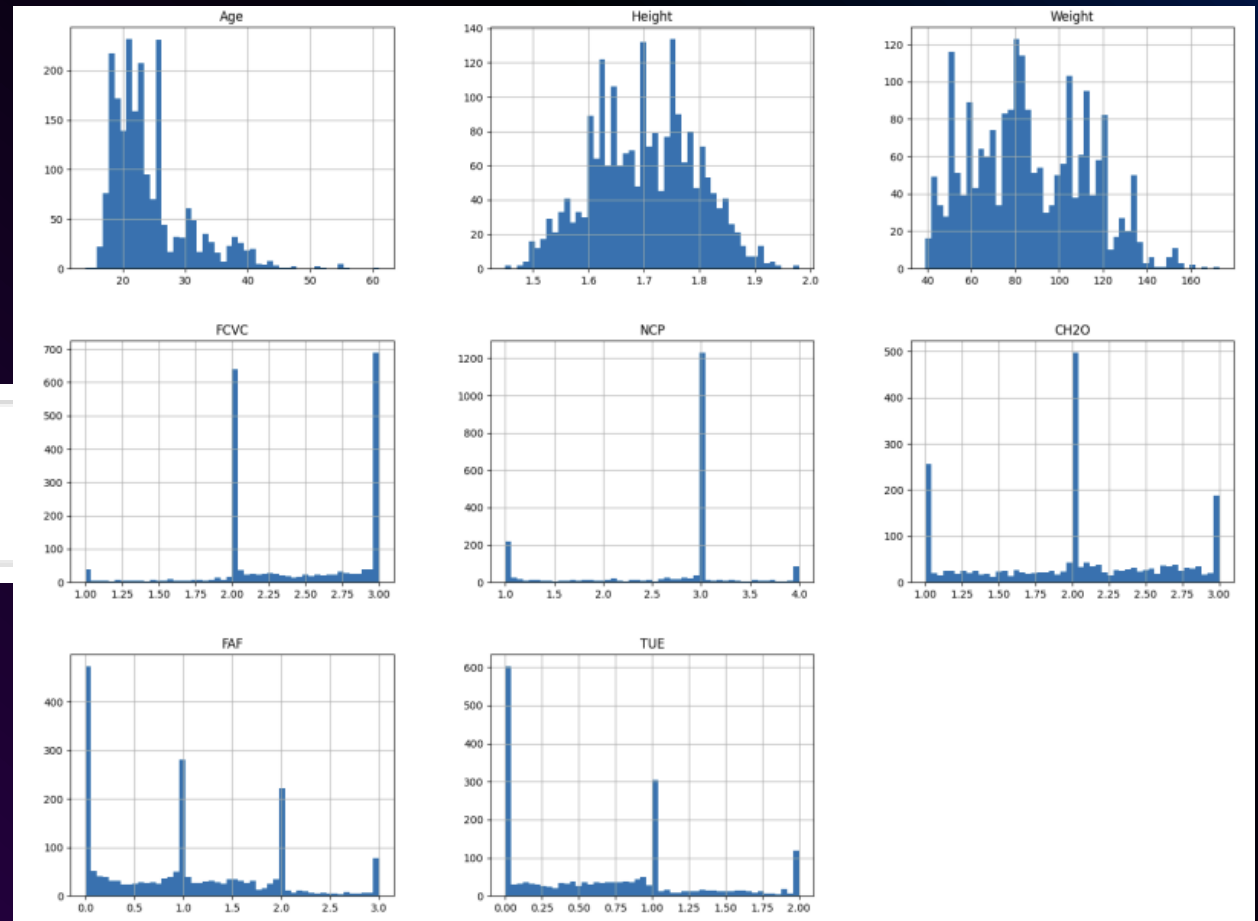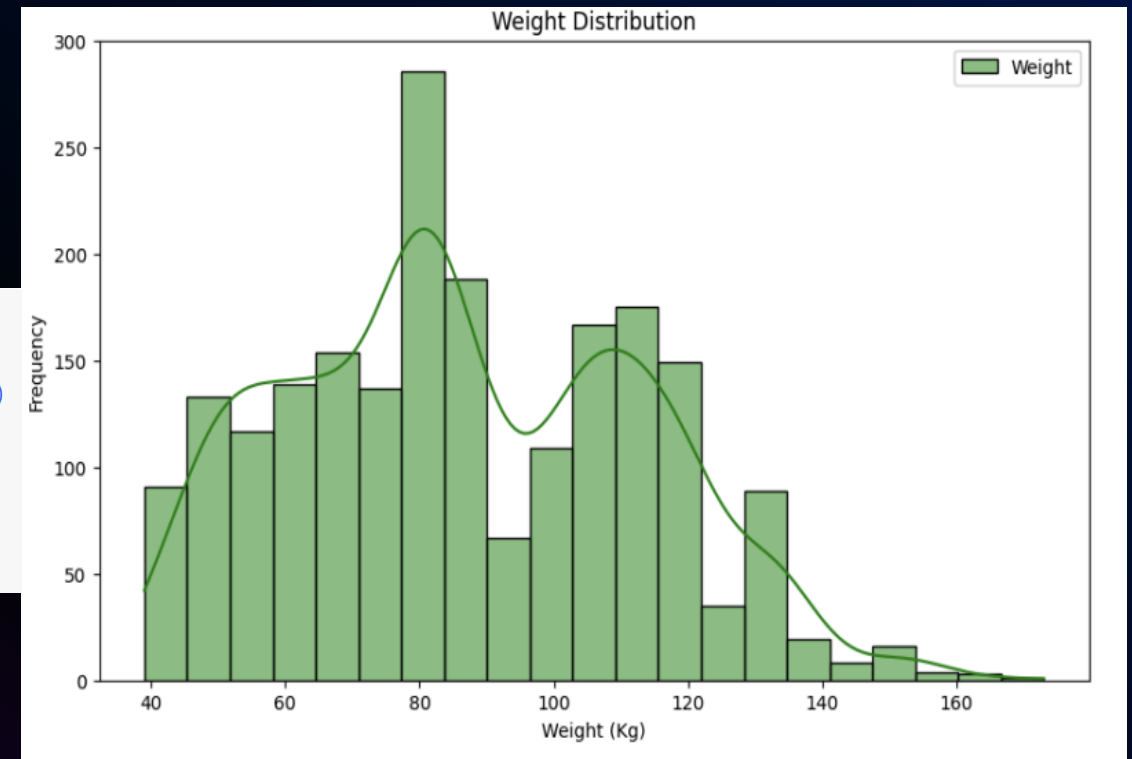
# GENERATING HISTOGRAM PLOTS FOR THE DATASET ATTRIBUTES



```
# Plot histograms for different parameters
data.hist(bins=50, figsize=(20,15));
```
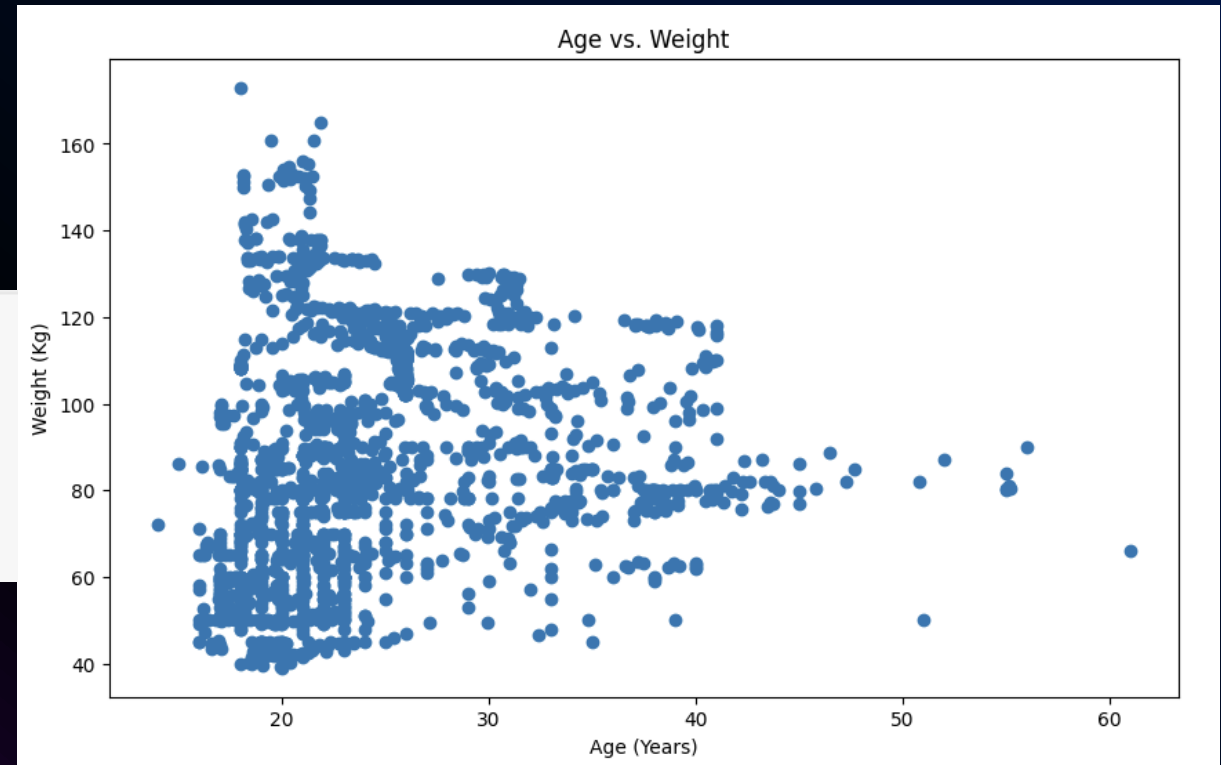
# HISTOGRAM PLOT SHOWING THE WEIGHT DISTRIBUTION

```python
1 # Plot Weight Distribution
2 plt.figure(figsize=(10, 6))
3 sns.histplot(data=data, x='Weight', color='green', kde=True, label='Weight')
4 plt.title('Weight Distribution')
5 plt.xlabel('Weight (Kg)')
6 plt.ylabel('Frequency')
7 plt.legend()
8 plt.show()
```
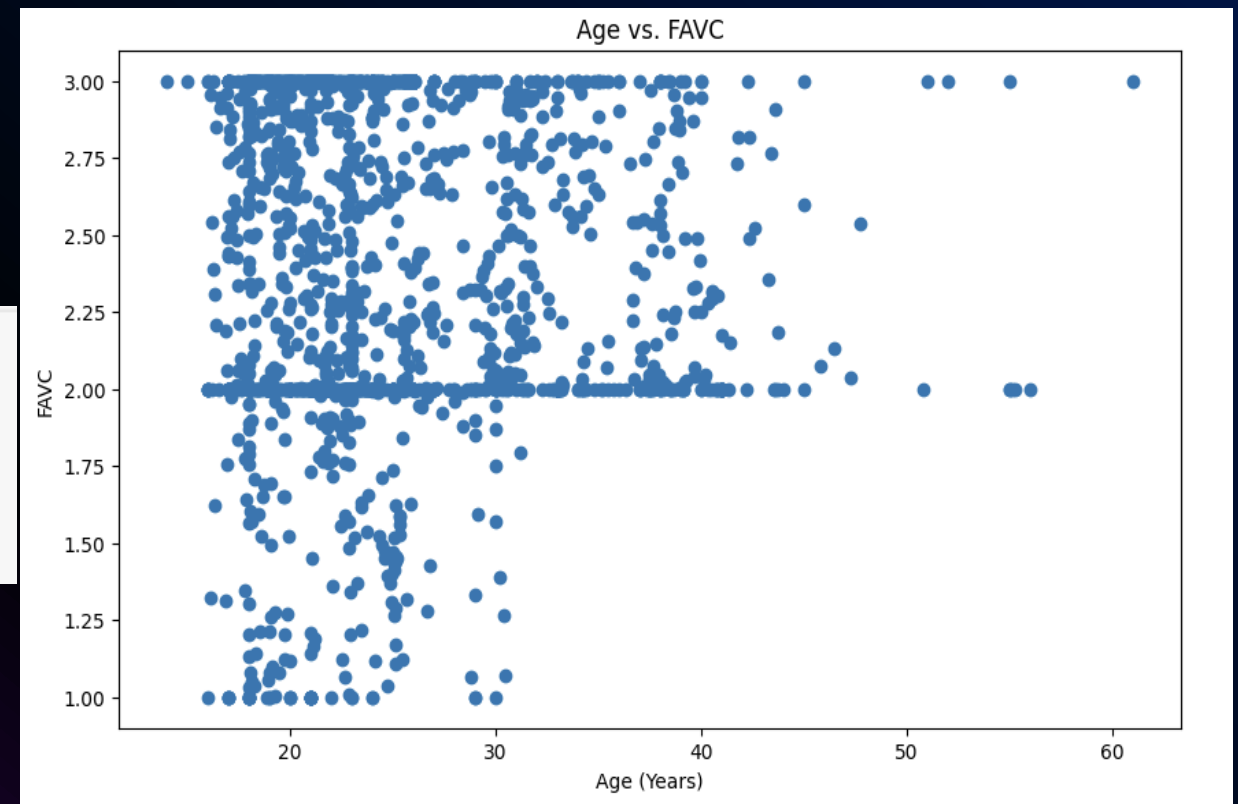
# SCATTER PLOT DEPICTING AGE VS. WEIGHT

```python
1 fig, ax = plt.subplots(figsize=(10,6))
2 ax.scatter(data['Age'], data['Weight'])
3 plt.title('Age vs. Weight')
4 ax.set_xlabel('Age (Years)')
5 ax.set_ylabel('Weight (Kg)')
6 plt.show()
```



Age vs. Weight

# SCATTER PLOT SHOWING AGE VS. FCVC

```python
1 fig, ax = plt.subplots(figsize=(10,6))
2 ax.scatter(data['Age'], data['FCVC'])
3 plt.title('Age vs. FAVC')
4 ax.set_xlabel('Age (Years)')
5 ax.set_ylabel('FAVC')
6 plt.show()
```



Age vs. FAVC

# SCATTER PLOT ILLUSTRATING AGE VS. HEIGHT FOR DIFFERENT WEIGHTS

# USING 'JOINTPLOT' FOR PLOTTING HEIGHT VS. WEIGHT

```
1 sns.jointplot(x="Height",y='Weight',data=data,kind='reg', color='green')
2 plt.title("Height vs Weight",loc='left')
```

# BAR PLOT SHOWING AGE VS. AVERAGE WEIGHT FOR THE TOP 15 AGES



Top 15 Ages with Highest Weight

```
1 # Top 15 Ages with Highest Weight
2 top_15_ages = data.groupby('Age')['Weight'].mean().nlargest(15)
3 plt.figure(figsize=(10, 6))
4 sns.barplot(x=top_15_ages.index, y=top_15_ages.values, color='blue')
5 plt.title('Top 15 Ages with Highest Weight')
6 plt.xlabel('Age (Years)')
7 plt.ylabel('Average Weight (Kg)')
8 plt.xticks(rotation=45)
9 plt.show()
```

```
1 data['Age']=data['Age'].round(2)
```

# BARPLOT ILLUSTRATING CONSUMPTION OF ALCOHOL FOR FEMALES & MALES

```python
1 data['CALC'] = data['CALC'].astype('category')
2 sns.barplot(x='CALC', y='Age',data=data,hue='Gender', palette='Set1')
3 plt.title('Consumption of Alcohol vs. Age')
4 plt.show()
```

# PIE CHART SHOWING THE % CONSUMPTION OF ALCOHOL

```
1 count = data['CALC'].value_counts()
2 labels = ["Always", "Frequently", "Sometimes", "No"]
3 vals = count.values
4 plt.pie(vals, labels=labels, autopct="%1.1f%%")
5 plt.title("% Consumption of Alcohol for Males & Females")
6
7 plt.show()
```



% Consumption of Alcohol for Males & Females

# PLOTTING AGE VS. NOBEYESAD USING SCATTER PLOT

```python
1 # Age vs. NObeyesdad
2 fig, ax = plt.subplots(figsize=(10,6))
3 ax.scatter(data['Age'], data['NObeyesdad'])
4 plt.title('Age vs. NObeyesdad')
5 ax.set_xlabel('Age (Years)')
6 ax.set_ylabel('NObeyesdad')
7 plt.show()
```



Age vs. NObeyesdad

# PLOTTING THE DISTRIBUTION OF CH2O

```python
print(data['CH2O'].describe())
plt.figure(figsize=(9, 8))
sns.distplot(data['CH2O'], color='g', bins=100, hist_kws={'alpha': 0.4});
```

```
count    2087.000000
mean        2.004749
std         0.608284
min         1.000000
25%         1.590922
50%         2.000000
75%         2.466193
max         3.000000
Name: CH2O, dtype: float64
```



https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

27

# COUNT PLOT SHOWING THE CATEGORICAL FACTOR 'CAEC'

# COUNT PLOT DEPICTING THE CATEGORICAL FACTOR 'CALC'

# COUNT PLOT DEPICTING NOBEYESDAD VS. FREQUENCY

# HOW TO INTERPRETING A BOX PLOT?



## Interpreting a box and whiskers

Construction of a box plot is based around a dataset's quartiles, or the values that divide the dataset into equal fourths. The first quartile (Q1) is greater than 25% of the data and less than the other 75%. The second quartile (Q2) sits in the middle, dividing the data in half. Q2 is also known as the median. The third quartile (Q3) is larger than 75% of the data, and smaller than the remaining 25%. In a box and whiskers plot, the ends of the box and its center line mark the locations of these three quartiles.

(25%) (25%) (25%) (25%)

Q1 Q2 Q3

IQR = Q3 - Q1    Upper whisker limit = Q3 + 1.5(IQR)    (Outliers)

# INSPECTING THE DATASET FOR OUTLIERS USING BOXPLOT

# INSPECTING THE DATASET FOR OUTLIERS USING BOXPLOT

# BOXPLOT SHOWING THE NCP FACTOR VALUES

# THE DATA SET FACTORS BEFORE REMOVING THE OUTLIERS



Example showing the removal of outliers using the diabetes.csv dataset
https://www.kaggle.com/datasets/saurabh00007/diabetescsv

# THE DATASET FACTORS AFTER REMOVING THE OUTLIERS



https://www.geeksforgeeks.org/detect-and-remove-the-outliers-using-python/

# WHAT IS A CORRELATION?

Correlation is a statistical indicator that quantifies the degree to which two variables change in relation to each other. It indicates the strength and direction of the linear relationship between two variables. The correlation coefficient is denoted by "r", and it ranges from -1 to 1.

- If r = -1, it means that there is a perfect negative correlation.
- If r = 0, it means that there is no correlation between the two variables.
- If r = 1, it means that there is a perfect positive correlation.

There are two popular methods used to find the correlation coefficients:

**Pearson's product-moment correlation coefficient**
The Pearson correlation coefficient (r) is a measure of linear relationship between two variables.

$$r = n(\sum xy) - (\sum x)(\sum y) / \sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}$$

Here,

- n is the number of data points
- $\sum xy$ is the sum of the product of corresponding values of x and y
- $\sum x$ is the sum of all the values of x
- $\sum y$ is the sum of all the values of y
- $\sum x^2$ is the sum of the squares of all values of x
- $\sum y^2$ is the sum of the squares of all the of y

# WHAT IS A CORRELATION MATRIX?

A correlation is a tabular representation that displays correlation coefficients, indicating the strength and direction of relationships between variables in a dataset. Within this matrix, each cell signifies the correlation between two specific variables. This tool serves multiple purposes, serving as a summary of data relationships, input for more sophisticated analyses, and a diagnostic aid for advanced analytical procedures. By presenting a comprehensive overview of inter-variable correlations, the matrix becomes invaluable in discerning patterns, guiding further analyses, and identifying potential areas of interest or concern in the dataset. Its applications extend beyond mere summary statistics, positioning it as a fundamental component in the preliminary stages of diverse and intricate data analyses.

https://www.geeksforgeeks.org/create-a-correlation-matrix-using-python/

# INTERPRETING THE CORRELATION MATRIX RESULTS

Strong correlations, indicated by values close to 1 or -1, suggest a robust connection, while weak correlations, near 0, imply a less pronounced association. They are identifying these degrees of correlation aids in understanding the intensity of interactions within the dataset, facilitating targeted analysis and decision-making. Positive correlations (values > 0) signify that as one variable increases, the other tends to increase as well. Conversely, negative correlations (values < 0) imply an inverse relationship—when one variable increases, the other tends to decrease. Investigating these directional associations provides insights into how variables influence each other, crucial for formulating informed hypotheses and predictions.

https://www.geeksforgeeks.org/create-a-correlation-matrix-using-python/

# CALCULATING THE PAIRWISE CORRELATION FOR ALL COLUMNS



https://www.geeksforgeeks.org/python-pandas-dataframe-corr/

# PLOTTING THE CORRELATION MATRIX HEATMAP

# CORRELATION MATRIX CLUSTER MAP



```
1 numeric_df = data.select_dtypes(include=['number'])# Select numeric columns only
2 # Calculate correlation matrix
3 correlation_matrix = numeric_df.corr()
4 sns.clustermap(correlation_matrix, annot=True, fmt=".2f")
5 plt.title("Correlation Between Features")
```

# CONVERTING CATEGORICAL VARIABLES TO NUMERIC VALUES

| | Gender | Age | Height | Weight | family_history_with_overweight | FAVC |
|---|---|---|---|---|---|---|
| 0 | Female | 21.000000 | 1.620000 | 64.000000 | yes | no |
| 1 | Female | 21.000000 | 1.520000 | 56.000000 | yes | no |
| 2 | Male | 23.000000 | 1.800000 | 77.000000 | yes | no |
| 3 | Male | 27.000000 | 1.800000 | 87.000000 | no | no |
| 4 | Male | 22.000000 | 1.780000 | 89.800000 | no | no |

| | Gender | Age | Height | Weight | family_history_with_overweight | FAVC | FCVC | NCP |
|---|---|---|---|---|---|---|---|---|
| 0 | Female | 21.0 | 1.62 | 64.0 | 1 | 0 | 2.0 | 3.0 |
| 1 | Female | 21.0 | 1.52 | 56.0 | 1 | 0 | 3.0 | 3.0 |
| 2 | Male | 23.0 | 1.80 | 77.0 | 1 | 0 | 2.0 | 3.0 |
| 3 | Male | 27.0 | 1.80 | 87.0 | 0 | 0 | 3.0 | 3.0 |
| 4 | Male | 22.0 | 1.78 | 89.8 | 0 | 0 | 2.0 | 1.0 |

Before

After

```
1 # replacing values
2 data['family_history_with_overweight'].replace(['no', 'yes'],
3                                   [0, 1], inplace=True)
4 data.head(5)
```

Converting categorical variables to numerical values for use in machine learning predictive models

https://www.geeksforgeeks.org/how-to-convert-categorical-variable-to-numeric-in-pandas/

# CONCLUSIONS
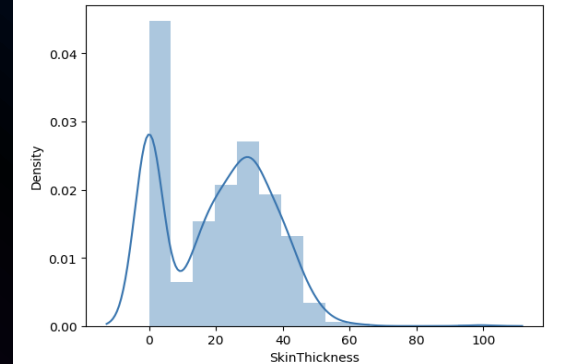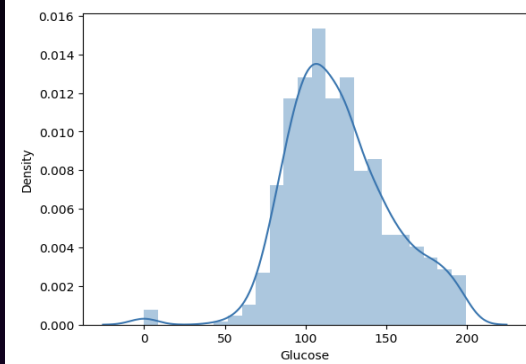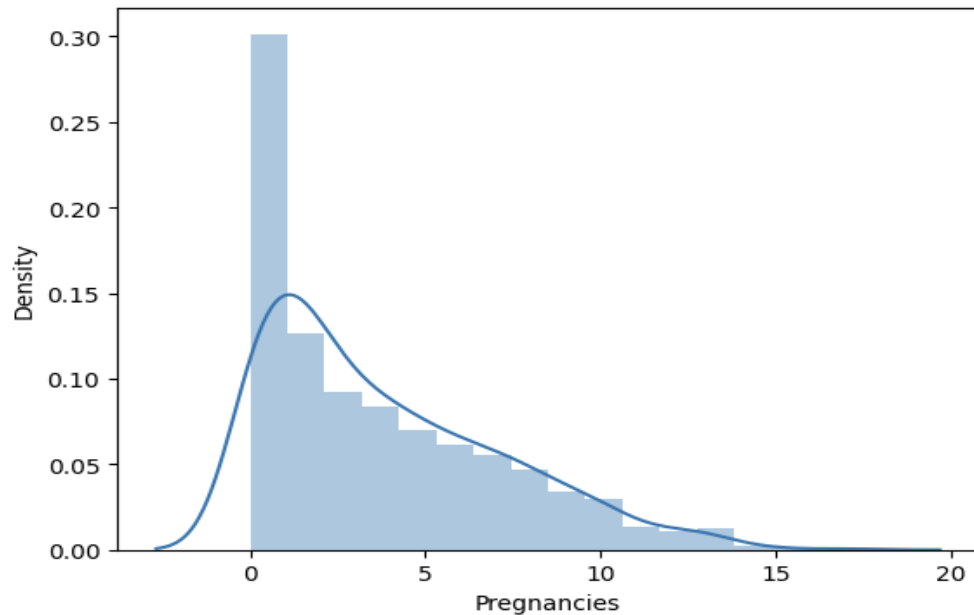
- The original dataset contains 2111 rows and 17 column
- There are 24 duplicate rows in the dataset
- All values in the dataset are unique and contain no null, NaN, **+** $\infty$ or missing values
- The NCP (number of main meals per day) data contains some outlier values and required removal
- There is a significant correlation between weight and height

# ADDITIONAL EXAMPLES

# DISTRIBUTION PLOT ILLUSTRATION

# PROBABILITY PLOT FOR TWO FACTORS IN THE DATASET

```python
1 import scipy
2 from scipy import stats
3 for feature in data.columns:
4     stats.probplot(data[feature], plot = plt)
5     plt.title(feature)
6     plt.show()
```
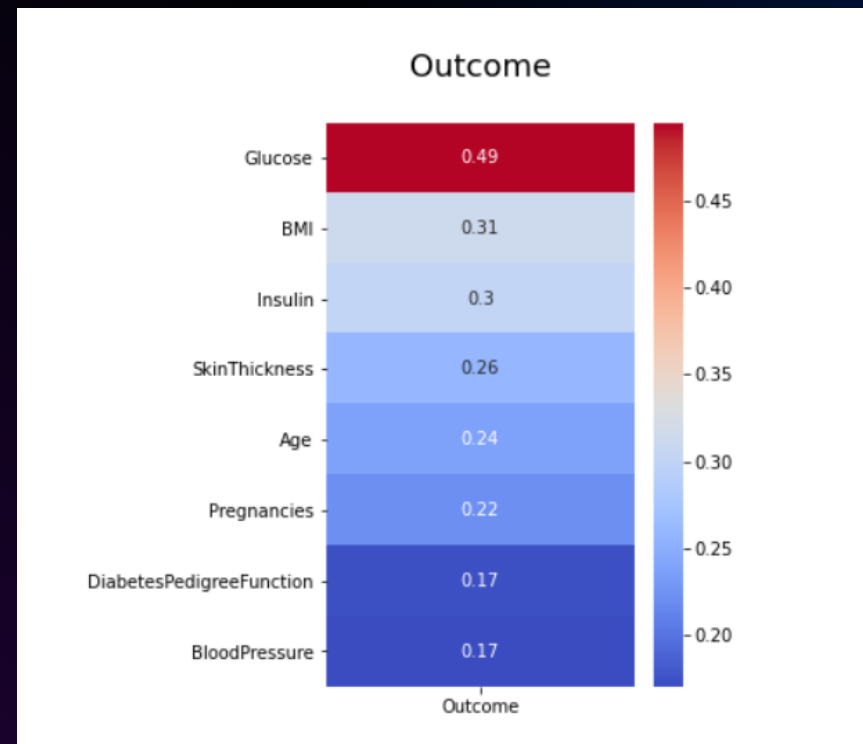


https://www.kaggle.com/datasets/mathchi/diabetes-data-set

47

# ANOTHER FORM OF HEATMAP PRESENTATION



```python
def corr_to_target(dataframe, target, title=None, file=None):
    plt.figure(figsize=(4,6))
    sns.heatmap(dataframe.corr()[[target]].sort_values(target,
                                                       ascending=False)
[1:],
                                                       annot=True,
                                                       cmap='coolwarm')

    plt.title(f'\n{title}\n', fontsize=18)

    plt.show();

    return

corr_to_target(df, "Outcome", title="Outcome")
```

https://www.kaggle.com/code/busekseolu/diabetes-classification

48

# REFERENCES

1.  **Data Science Horizon, Data cleaning and preprocessing for data science beginners.** (https://datasciencehorizons.com/data-cleaning-preprocessing-data-science-beginners-ebook/)

2.  **Matthieu Komorowski, et al., Exploratory Data Analysis, Chapter 15, doi:10.1007/978-3- 319-43742-2_15.**

3.  **DataQuest, Data Science Cheat Sheet-Pandas.** (https://s3.amazonaws.com/dq-blog-files/pandas-cheat-sheet.pdf)

4.  **DataCamp, Python for Data Science Cheat Sheet-Matplotlib.** (https://s3.amazonaws.com/assets.datacamp.com/blog_assets/Python_Matplotlib_Cheat_Sheet.pdf)

5.  **DataQuest, Data Science Cheat Sheet-Numpy.** (https://assets.datacamp.com/blog_assets/Numpy_Python_Cheat_Sheet.pdf)

6.  **DataCamp, Python for Data Science, Seaborn Cheat Sheet.** (https://s3.amazonaws.com/assets.datacamp.com/blog_assets/Python_Seaborn_Cheat_Sheet.pdf)

7.  **DataCamp, Python for Data Science Cheat Sheet-Scikit-Learn.** (https://s3.amazonaws.com/assets.datacamp.com/blog_assets/Scikit_Learn_Cheat_Sheet_Python.pdf)

8.  **David Beazley, et al., Python Cookbook, 3rd Ed., O'Reilly, Beiing, 2013.**